

# 01TXFSM - Machine Learning and Deep Learning

## Final Project First Person Action Recognition

Eros Fani - s269781  
Politecnico di Torino

eros.fani@studenti.polito.it

Gabriele Trivigno - s276807  
Politecnico di Torino

gabriele.trivigno@studenti.polito.it

Cristiano Gerbino - s277058  
Politecnico di Torino

s277058@studenti.polito.it

### Abstract

## 1. Introduction

### 1.1. Goals

The first goal of the project is to replicate some of the experiments performed in [3] and [2]. The objective of these studies is the First Person Action Recognition: they tried to implement a deep learning model capable to extract meaningful features to automatically predict the action filmed by a wearable camera.

After having replicated these experiments we performed a grid search on the experiments to find the best set of values for the hyperparameters.

At last we have tried to improve the performances of the results of [3] and [2] with some innovative ideas.

### 1.2. Our contribution

...

### 1.3. Data exploration

The dataset under analysis is a modified version of GTEA61<sup>1</sup>. The dataset contains the videos in form of frames, and also two kind of preprocessed images: *motion maps* and *optical flows*. The folder schema of the dataset is shown in Figure 1. Videos represent 61 class actions performed by 4 different users (*S1*, *S2*, *S3*, *S4*). Sometimes for some actions more than one video is available. The total number of videos in the dataset is, however, 457, which actually means that it is a quite small dataset.

The optical flow methods try to calculate the motion between two image frames which are taken at times  $t$  and  $t + \Delta t$  at every voxel position. The warp flow methods try also to remove the motion of the wearable camera. We have two kind of these last representations in our dataset: one computed in the horizontal axis (folder *flow\_x\_processed*) and one other computed in the vertical axis (folder *flow\_y\_processed*).

The motion maps are special black-and-white images which represent the spatial location in which the Motion Segmentation task of [2] focuses its attention per each frame. The mmaps present large similarities with the warp flows.

<sup>1</sup>Georgia Tech Egocentric Activity Datasets: <http://cbs.ic.gatech.edu/fpv/>

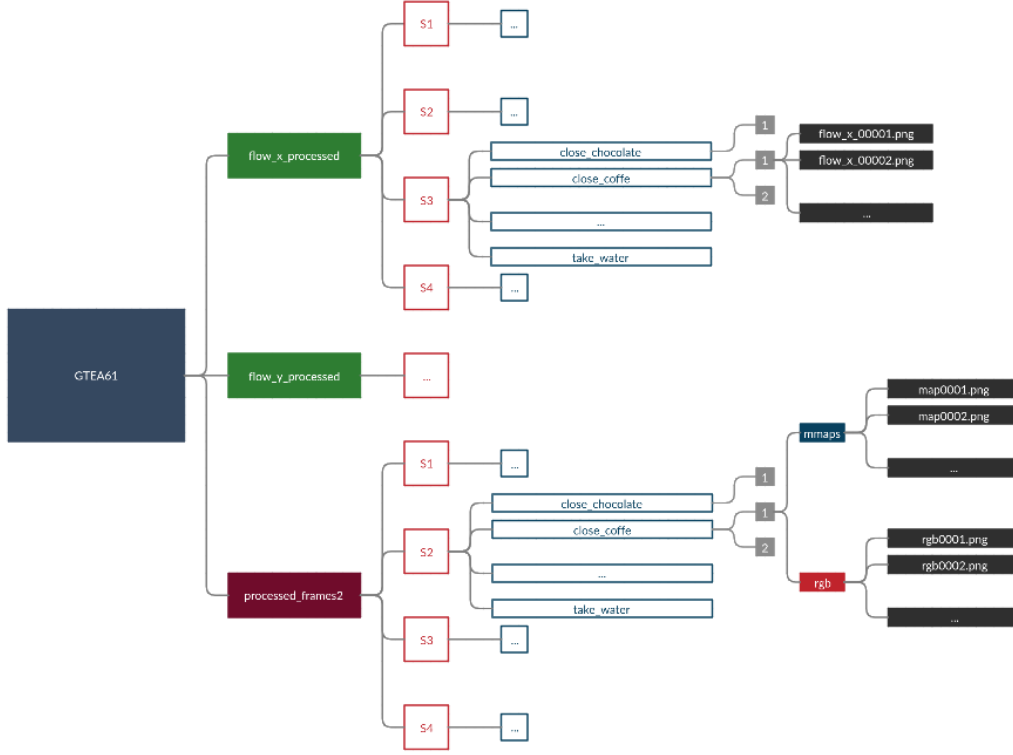


Figure 1: Folder schema of our GTEA61

The differences between the kind of available images in our dataset are shown in Figure 2.

## 1.4. Data cleaning

The dataset was almost clean already from the beginning, but we encountered two problems within it:

- there were hidden useless folders *.DSstore* inside each one of the user folders. These have been removed
- some of the first mmmaps of some videos were missing. In these cases we have simply duplicated the second mmap

## 2. Descriptions of the models

Here we describe the models that we have used to perform our experiments.

### 2.1. Egornn

Egornn is a Recurrent Neural Network. The overall architecture of *egornn* is shown in Figure 3. This net is based on *resnet34*[1], which constitutes the main block. *resnet34* has five convolutional layers inside itself: with respect to Figure 3 they are: *Conv*, *Layer1*, *Layer2*, *Layer3* and , *Layer4*. From now on we'll call these blocks respectively *conv1*, *conv2*, *conv3*, *conv4* and *conv5*.

At the termination of the *resnet34* is placed a *Spatial Attention Layer*. It includes a *Class Activation Map* (CAM) that is capable to identify the image regions that have been used by the CNN to identify



Figure 2: Types of images in our dataset. In this example is shown a sample of images from the *close\_chocolate* action. From the left column to the right column: rgbs, warp flows x, warp flows y, motion maps

the class under analysis. It is computed by taking the output of the *softmax* layer and the output of *conv5* and taking the linear combination of all the weights of *conv5* and the weights of the softmax.

The output of the CAM is then sent to a *softmax* layer to obtain a probability map, which is called *Spatial Attention Layer* (SAM). The output of the SAM is finally multiplied, cell by cell (Hadamard product), with the output of *conv5*, obtaining another tensor of weights which is sent to a *Convolutional Long Term Support Memory* block (ConvLSTM).

The reason for the usage of the ConvLSTM block is that, up to now, what the net does is to take each frame and to try to make predictions based only on the features that the net can extract from those frames, without taking into consideration the temporal encoding of frame level features. The convLSTM block take into consideration, for each frame  $i$ , both the output of the SAM for the layer  $i$  and the output of the ConvLSTM for the layer  $i - 1$ , constituting a recursive structure.

The last output of the ConvLSTM (the output obtained from the last frame of a particular video) is average pooled and reshaped to obtain a final classification layer with 61 neurons (i.e. the number of classes of our dataset).

## 2.2. Flow\_resnet34

*Flow\_resnet34* is just a *resnet34* edited to work with the warp flows. It gets five warp flows from *processed\_frames\_x* and five from *processed\_frames\_y* in form of a tensor of ten channels and tries to make predictions on the 61 classes.

## 2.3. Two stream model

Egornn learns appearance features, while *flow\_resnet34* learns motion features. The way to join the two nets is to concatenate the two output layers and to add at the end a fully connected layer to get the class category scores.

## 2.4. Motion Segmentation branch applied to egornn

The problem which [2] tries to overcome is that in the two stream model motion and appearance are actually separately learned, without taking into account the spatial-temporal relationships.

We have built an architecture similar to *sparnet*, where the *motion segmentation block* is the same but the *action recognition block* has been substituted by egornn (like in one of the attempts in [2]). The architecture is shown in Figure 4. We have

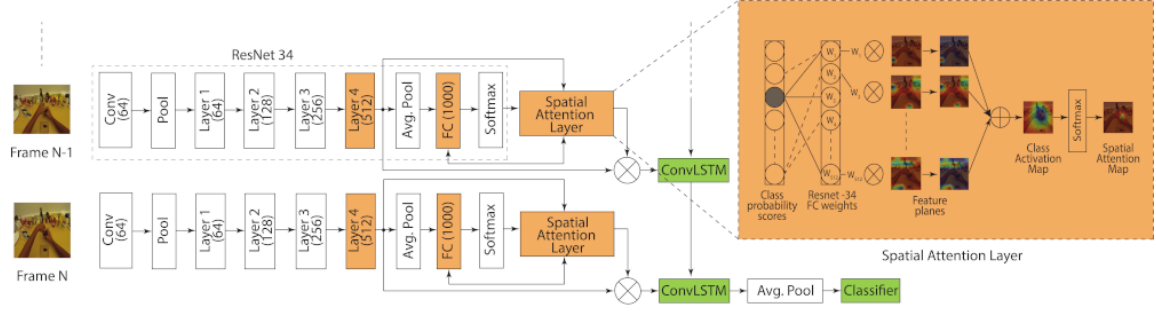


Figure 3: Architecture of *egornn*

used this architecture with some granular variations during our experiments, but the main blocks are always as shown in Figure 4. The input of the convolutional layer of MS Block is taken from one of the convolutional layers of *resnet34* of *egornn* (the actual layer varies with our experiments). Then, after the convolutional layer, there is a fully connected layer followed by a softmax which normalizes the weights between 0 and 1. *mmaps ground truth* and *rgb ground truth* represent the mmaps and the rgb after the transformations. The transformations applied to the mmaps are the same applied to the rgbs, plus a small amount of proper mmaps transformations which always ends with a transformation which linearizes the pixels (from a 2 dimensional tensor per mmap we get a 1 dimensional tensor per mmap). For the msblock these linearized pixels represent a real ground truth, because each of the output neurons of the MS Block is used to predict the values of the mmaps ground truth. The pixel losses are summed together (obtaining as result  $L_{ms}$ ) and then are summed again with the *egornn* loss ( $L_c$ ). The final loss is used to compute the gradients to update the weights.

### 2.5. Static-dynamic discriminator

Starting from the model described above, we have added a final binary classifier to *egornn* after the convLSTM, parallel to the other classifier already present which still tries to predict the actual class of the video. The idea was to force the net to learn the motion features from the rgbs. During the training phase this classifier gets 2 kind of sequences of frames: one is the same of the original classifier, while the other gets a sequence of identical frames. This discriminator should be able to recognize the

actual videos from the static frames. In this way the gradients should adapt to focus the attention on the motion.

## 3. Experiments

Our nets are always trained on a predefined train set, which includes all and only the videos of the users *S1*, *S3* and *S4*, while validation and test sets coincide and is constituted by all and only the videos of a single user, *S2*. In addition, the weights of the *resnet34* are pretrained on ImageNet.

Due to Colab limitations of GPU memory, we have only been able to perform experiments on a limited amount of frames (7 or, in less cases, 16). Due to this problem our results should be interpreted not as absolute value of the accuracy, but as a sort of relative value with respect to the number of frames for each video in our batches.

The size of our batches has always been left to 32, as well as the number of hidden units of the convLSTM module, fixed at 512. Our optimization algorithm is always Adaptive moment estimation (ADAM) with the only exception of *flow\_resnet34*, for which it is Stochastic Gradient Descent (SGD). When using this last optimizer, the momentum has always been left to 0.9. The scheduler is a MultiStepLR scheduler, which decreases the original learning rate LR by a factor GAMMA at each value of STEP\_SIZE.

### 3.1. Egornn

We have replied some of the same experiments of [3] on the original *egornn*. We have run each of

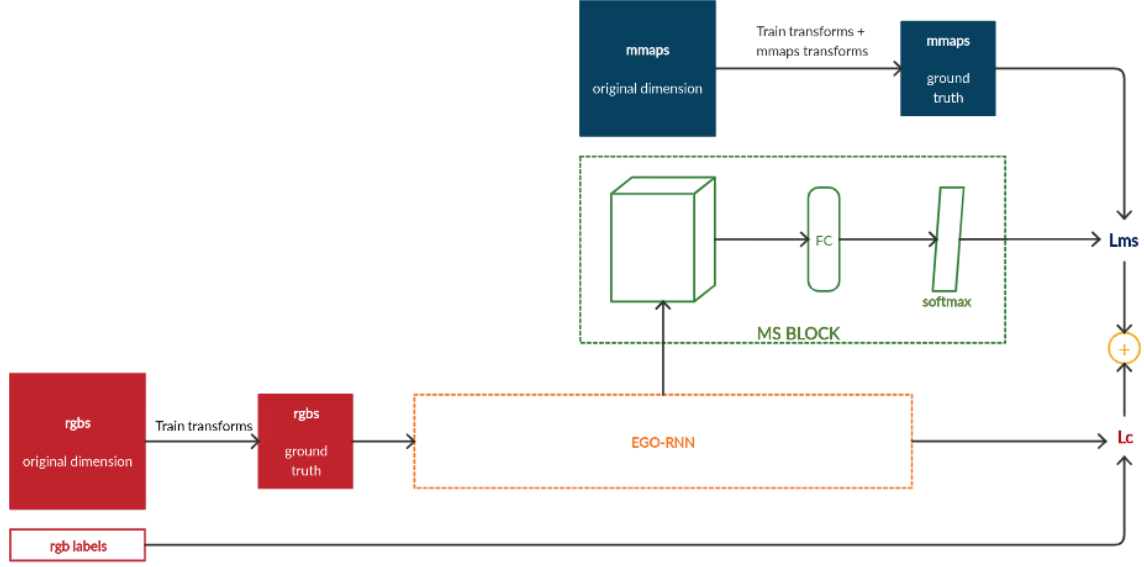


Figure 4: Generic architecture of motion segmentation branch applied to *egornn*

these experiments three times and then we have averaged the results.

First, we have performed the classification by using the *egornn* without and with the CAM. The training phase has been divided in two parts, as in the original paper:

1. train of ConvLSTM and Classifier (green blocks in Figure 3)
2. train of conv5 (layer4 of *resnet34*), FC(1000), Spatial Attention Layer (orange blocks in Figure 3) in addition to the previously listed blocks

The values of the hyperparameters for the first stage are:

<b>LR</b>	1e-3
<b>WEIGHT_DECAY</b>	4e-5
<b>NUM_EPOCHS</b>	200
<b>STEP_SIZE</b>	[25, 75, 150]
<b>GAMMA</b>	0.1

While, for the second stage, they are:

<b>LR</b>	1e-4
<b>WEIGHT_DECAY</b>	4e-5
<b>NUM_EPOCHS</b>	200
<b>STEP_SIZE</b>	[25, 75]
<b>GAMMA</b>	0.1

Then, we have also trained *flow\_resnet34* alone. In this case we used only 5 frames per each flow (x and y) due to the fact that for some videos no more than 5 frames were provided.

At last we performed the two stream training.

The summary of our results is shown in Figure 5.

### 3.2. Motion Segmentation branch applied to *egornn*

## References

- [1] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition, 2015.
- [2] M. Planamente, A. Bottino, and B. Caputo. Joint encoding of appearance and motion features with self-supervision for first person action recognition, 2020.
- [3] S. Sudhakaran and O. Lanz. Attention is all we need: Nailing down object-centric attention for egocentric activity recognition, 2018.

Configurations	Frames	Mean accuracy
EGO-RNN without CAM - stage 1	7	29.89
EGO-RNN without CAM - stage 1	16	27.87
EGO-RNN without CAM - stage 2	7	50.00
EGO-RNN without CAM - stage 2	16	50.57
EGO-RNN - stage 1	7	41.38
EGO-RNN - stage 1	16	46.84
EGO-RNN - stage 2	7	58.91
EGO-RNN - stage 2	16	65.52
flow_resnet34	5	49.71
two-stream (joint train)	7*	57.76
two-stream (joint train)	16*	66.38

Figure 5: Summary of the results over different configurations. Each value of the mean accuracy is the mean of the accuracies over three identical experiments. \*the number of frames refers to the egornn branch (for the flow\_resnet34 branch the number of frames is always 5)